

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-108382

(43)Date of publication of application : 10.04.2002

(51)Int.Cl.

G10L 15/00  
G06T 13/00  
G06T 15/70  
G10K 15/02  
G10L 13/00  
G10L 21/06  
G10L 15/22

(21)Application number : 2000-294151

(71)Applicant : SONY CORP

(22)Date of filing : 27.09.2000

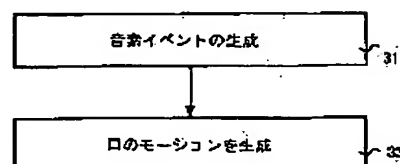
(72)Inventor : OTO YASUNORI  
UEDA YUICHI

(54) ANIMATION METHOD AND DEVICE FOR PERFORMING LIP SYNCHRONIZATION

(57)Abstract:

PROBLEM TO BE SOLVED: To provide animation technique capable of automatically generating an animation automatically including a mouth part according to voice data, etc.

SOLUTION: Voice data are received and a phoneme analysis is taken to generate phoneme events along a time base (31). A mouth-shape animation corresponding to one or more phonemes is weighted and added to generate an animation of a distinctive or vague mouth shape corresponding to an articulate or obscure phoneme (32). For a phoneme whose mouth movement is earlier than pronunciation, the phoneme event is made earlier than voice data timing.



## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号  
特開2002-108382  
(P2002-108382A)

(43) 公開日 平成14年4月10日 (2002. 4. 10)

(51) Int.Cl. <sup>7</sup>	識別記号	F I	テマコード (参考)
G 1 0 L 15/00		G 0 6 T 13/00	B 5 B 0 5 0
G 0 6 T 13/00		15/70	B 5 D 0 1 5
15/70		G 1 0 K 15/02	5 D 0 4 5
G 1 0 K 15/02		G 1 0 L 3/00	5 5 1 H
G 1 0 L 13/00			Q
審査請求 未請求 請求項の数11 O L (全 9 頁) 最終頁に続く			

(21) 出願番号 特願2000-294151 (P2000-294151)

(22) 出願日 平成12年9月27日 (2000. 9. 27)

(71) 出願人 000002185

ソニー株式会社

東京都品川区北品川6丁目7番35号

(72) 発明者 大戸 康紀

東京都品川区東五反田1丁目14番10号 株式会社ソニー木原研究所内

(72) 発明者 上田 裕一

東京都品川区東五反田1丁目14番10号 株式会社ソニー木原研究所内

(74) 代理人 100101801

弁理士 山田 英治 (外2名)

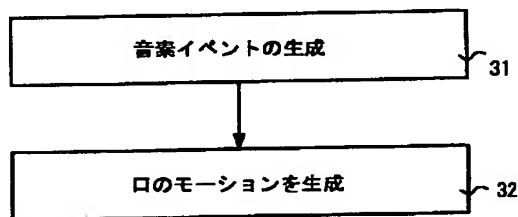
最終頁に続く

(54) 【発明の名称】 リップシンクを行うアニメーション方法および装置

(57) 【要約】

【課題】 リップシンク・アニメーションを簡易に生成する。

【解決手段】 音声データを受けとって音素解析を行い、時間軸に沿って音素イベントを生成する (31)。1または複数の音素に対応した口形状アニメーションを重み付け加算して明瞭な、あるいは、あいまいな音素に対応した明瞭な口形状、あるいは、あいまいな口形状のアニメーションを生成する (32)。口の動きのほうが発声より早い音素の場合には音素イベントを音声データタイミングより早めにする。



**【特許請求の範囲】**

**【請求項1】** 音声データから音素データを生成するステップと、

上記音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成するステップとを有することを特徴とするアニメーション生成方法。

**【請求項2】** 上記口の部分のアニメーションを生成するステップは、音素データの組合せによりあいまいな音声を表し、この組合せに含まれる音素データにそれぞれ対応する口の部分のアニメーションを重み付け合成する請求項1記載のアニメーション生成方法。

**【請求項3】** 所定の音素データについては、音素データのタイミングより、対応する口の部分のアニメーションの生成のタイミングが所定時間早くする請求項1記載のアニメーション生成方法。

**【請求項4】** 音素データの生成を有限状態オートマトンを参照して行い音素間の遷移に所定の制約を設ける請求項1記載のアニメーション生成方法。

**【請求項5】** 音声データから音素データを生成するステップと、

上記音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成するステップと、上記音素データを基準にして上記音声データと上記アニメーションとを同期させて再生するステップとを有することを特徴とするアニメーション生成方法。

**【請求項6】** 音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成し、さらに、音素データの組合せによりあいまいな音声を表し、この組合せに含まれる音素データにそれぞれ対応する口の部分のアニメーションを重み付け合成することを特徴とするアニメーション生成方法。

**【請求項7】** 音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成し、さらに、所定の音素データについては、音素データのタイミングより、対応する口の部分のアニメーションの生成のタイミングが所定時間早くすることを特徴とするアニメーション生成方法。

**【請求項8】** 音素データの組合せによりあいまいな音声を表し、この組合せに含まれる音素データにそれぞれ対応する口の部分のアニメーションを重み付け合成する請求項7記載のアニメーション生成方法。

**【請求項9】** 音声データから音素データを生成する手段と、

上記音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成する手段とを有することを特徴とするアニメーション生成装置。

**【請求項10】** 音声データから音素データを生成する手段と、

上記音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成する手段と、

上記音素データを基準にして上記音声データと上記アニメーションとを同期させて再生する手段とを有することを特徴とするアニメーション生成装置。

**【請求項11】** 音声データから音素データを生成するステップと、

上記音素データに基づいて二次元ないし三次元のモデルの口の部分のアニメーションを生成するステップと、上記音素データを基準にして上記音声データと上記アニメーションとを同期させて再生するステップとをコンピュータに実行させるために用いるコンピュータプログラムをコンピュータ読取り可能に記録するアニメーション生成用の記録媒体。

**【発明の詳細な説明】**

**【0001】**

**【発明の属する技術分野】** 本発明は、二次元あるいは三次元形状のモデルにおいて、音声データに同期して、口となる部分のアニメーションを生成する技術に関する。本発明の技術はアニメーションの作成時あるいは実行時に使用され、ゲームや映像コンテンツなどで活用される。

**【0002】**

**【従来の技術】** 従来のリップシンクアニメーションでは、音声データを記録した後、そのデータに合わせてリップアニメーションを作成する必要があった。もしくは、予めあたりをつけておいた口のアニメーションに対して音声データを当てていく必要があった。

**【0003】**

**【発明が解決しようとする課題】** 従来のリップシンクアニメーションでは、音声データを見ながらアニメーションを作成していく場合、かなりの入力が必要であった。また、この場合、音声データとアニメーションとを分離することが難しく、リップアニメーションの再利用性もほとんどなかった。また、音声データを追加する場合、新たにリップアニメーションを作成したり、前に作成したアニメーションを修正する必要があった。

**【0004】** また、アテレコによって音声データをアニメーションに付加していく場合、アニメーションと音声の間に明確な対応関係がなく、声優の経験によるところが大きかった。

**【0005】** 本発明は上述の事情を考慮してなされたものであり、音声データ等に基づいて自動的に口の部分を含むアニメーションを自動的に生成することができるアニメーション技術を提供することを目的としている。

**【0006】**

**【課題を解決するための手段】** 本発明によれば、上述の目的を達成するために、特許請求の範囲に記載のとおり構成を採用している。ここでは、本発明について若干補充的な説明を行う。

**【0007】** 本発明の原理的な構成例によれば、音声データの音素解析を行った後、解析で得られた音素の種類

に対する重みを計算し、あいまいな口形状を示すデータを作成する。次に、音素間の形状変化アニメーションの作成を行い、音声データに対してイベント作成した口形状アニメーションをマッピングする。この際、破裂音（／p／や／b／）などにおいて口形状の変化が音素の発音以前に行われるため、イベントを前方に修正する。口形状のアニメーションを行う時点で、音声データを再生し、これと同時に、イベントに基づいてアニメーション合成を行うことにより、同期の取れたアニメーションの表示を行う。

【0008】このような構成によれば、二次元あるいは三次元モデルのアニメーションの作成とは別に、音声データのみを後から独立して追加することが可能であり、モデル設計時においてリップシンクアニメーションを個々に作成する必要がなくなる。さらに、音声データに対して、インタラクティブにリップアニメーションを生成することができる。また、／a／と／e／の中間などのあいまいな発音に対して、あいまいな口の動きを実現することができる。また、子音に対して、有声音と無声音（／z／と／s／など）や、破裂音と摩擦音（／p／と／z／）など、発音状態として同時にその状態を取り得ない場合や、音素間の連続性として、／y, a／、／y, y／、／y, o／などが許されるのに対して、／y, i／、／y, e／が許されないなどといった制約を盛り込むことが可能になる。

【0009】なお、音素としては、子音、母音単位であつかつてよいし、より大きな単位例えば音節のシンボルとして扱ってもよい。要するに口形状とマッピング可能な単位のものであればどのような音素を用いてもよい。

【0010】また、本発明は方法としても装置としても実現可能である。また、そのような方法をコンピュータで実現するために用いるプログラムを記録したコンピュータ読取り可能な記録媒体も本発明の技術的な範囲に含まれる。

#### 【0011】

【発明の実施の形態】ここでは、本発明におけるリップシンクアニメーション生成およびその表示を実現するための装置の説明をした後、音声データの解析方法と、あいまいな発音に対する、口のあいまいなアニメーションの生成方法について説明を行う。

【0012】図1は、本発明の実施例のアニメーション生成装置を全体として示しており、この図において、アニメーション生成装置1は、音素イベント生成部（アプリケーション）2、アニメーション処理部（アプリケーション）3、オペレーティングシステム4、入力装置5、出力装置6、その他のリソース（ハードウェア、ソフトウェア）7等を含んで構成されている。アニメーション生成装置1は、実際には、ゲーム機、マルチメディアキオスク、パーソナルコンピュータ等を実装される。

アニメーション編集装置として構成してもよい。オペレーティングシステム4は、実装環境に応じたものであり、パーソナルコンピュータ用の汎用のオペレーティングシステムでもよいし、機器独自の組み込みオペレーティングシステムでもよい。音素イベント生成部2は、入力音声データを解析してアニメーション生成の音素イベントを生成するものである。詳細については図4を参照して後に説明する。アニメーション処理部3は、音素イベント生成部2から音素イベントを受けとってリップシンクのアニメーションを合成して画像データを生成するものである。図示しないが、画像生成の一部を、専用のハードウェアを用いて行ってもよい。

【0013】図2は、図1の音素イベントおよびアニメーション処理の関係を説明するものである。図1において、音声データを解析して音素データ21（「k o」、「n n」、「n i」、「t i」、「w a」）が生成され、音素データの出現タイミングに応じて対応するアニメーション22、23、24が生成される。各アニメーション22、23、24は音素に応じた口形状をしており、音素イベントの出現に応じて対応するアニメーションを起動することによりリップシンクを実現できる。

【0014】図3は、本実施例の動作を概要を示しており、この図に示すように、本実施例では、音声データに対する前処理を行って音素イベントを生成し（31）、この後、音声データに同期した口のアニメーションを行っている（32）。もちろん、音声データに対する前処理は逐次処理が可能であり、音声データをストリーミング入力しながら、アニメーションを生成することも可能となっている。これらの処理31および32は図1の音素イベント生成部2およびアニメーション処理部3にそれぞれ対応しており、それぞれ後に詳述する。

【0015】図4は、音声データに対して行う前処理の流れを示す。図4において、音声データをシステムに入力し（401）、音素解析を行った後（402）、状態遷移を通すことによってその候補を制限する（403）。同時に複数の音素が候補として残る場合には、これによって得られた複数の音素分に対して口形状の重み割合を計算し、口形状の合成情報を作成する（404）。この後、音素遷移におけるアニメーションを作成し（405）、音素データに対してイベントとアニメーションの登録を行う（406）。

【0016】次に、これらの処理について一つずつ説明を行う。

【0017】図5に音声データの入力ソース（506）の一例を示す。インターネットにおけるストリーミングデータ（501）や放送（502）、マイクによる直接入力（503）や、CD（504）、MD（505）などの記録媒体からの入力が可能となっている。

【0018】次に、音素解析について説明する。最初に図6に示すように音声の波形データ（602）におい

て、ゼロに交わる点(601)を抽出し、その周期性を調べる。子音が過渡的であるのに対し、母音部における周期性が揃っていることから、解析フレームを適宜作成していくことにより、母音・子音の分離と、スペクトル解析の精度を向上させることができる。ここで、603は波形の周期を示している。また604で示す部分は、波形の周期が過渡的であることから子音フレームとして判断されている。また、605で示す部分は波形の周期が揃っていることから、母音フレームと判断されている。

【0019】次に、図7に示すように、入力された音声データ(a)において、図6で判断した解析フレームごとに周波数スペクトル(b)、(c)を求める。

【0020】なお、周波数スペクトルの求め方としてはFET(高速フーリエ変換)やMEM(最大エントロピー法)などがある。これらの手法は周知であるのでとくに説明は行わない。詳細については例えば「時系列解析プログラム」(北川源四郎著、岩波書店発行)を参照されたい。

【0021】取得した周波数スペクトルはいくつかの主要となる周波数成分を持っており、図8(a)に示すように、低周波数側から、第一フォルマント(801)、第二フォルマント(802)、第三フォルマント(803)と言われている。本実施例ではこれらのフォルマント間の関係を用いて音素解析を行う。なお、フォルマントに関しては、「フーリエの冒険」(トランスナショナルカレッジオブブックス編集、ヒポファミリークラブ発行)等を参照されたい。

【0022】ここで、リップシンクアニメーションにおいては、同一モデルに対して一人の音声データ提供者が担当することが普通であり、この実施例では、特定話者を想定することにする。なお、不特定話者の場合に対しても、音素解析の手法が変化するだけであり、全体としては同一の構成となる。

【0023】そして、特定話者を想定することから、図8(b)に示すように、フォルマントと音素の関係についてキャリブレーションを行っておく。ここで、804~808は、各母音の位置を示している。また、これは話者の音質が変化しない限りにおいて、再度取り直す必要がない。

【0024】図9(a)に示すように、音声データに対して設定した解析フレーム毎のフォルマントデータ91を取得し、図8(b)において行った、キャリブレーションしたフォルマント位置との関係を取得する。

【0025】次に、図10に示すような音素間の状態遷移を考慮し、音素候補の絞り込みを行う。図10(a)の状態遷移では、1001は音素/k/を表し、すべての母音(図1002~1006)への変化が可能であることを示している。一方、図10(b)の状態遷移では、1007は音素/y/を表し、/a/(1008)、/

u/(1009)、/o/(1010)の母音へのみの状態遷移が許されていることを示している。

【0026】また、図11における、図1101から1102の遷移と、1103から1104への遷移によって生成されるアニメーション(1105から1106への変化)が同じものについては、一つにまとめて扱うことができる。

【0027】次に、図12に示すように、音素候補(1202, 1203, 1205, 1206, 1207)におけるフォルマント位置と計測されたフォルマント位置(図1204)の距離を計算し、各音素候補毎の重みを計算する。重み計算例を図1201に示す。この際、人間の耳は実際の周波数より、オクターブとして認識することから、各成分に対して対数を取った上で扱っている。

【0028】また、図13に示すように、第三フォルマントを用いる場合も、距離の計算方法(1301)が変化する以外は、音素候補(1302, 1303, 1305, 1306)におけるフォルマント位置と計測されたフォルマント位置(1304)から、同様に計算することができる。

【0029】次に、図14に示すように、各音素フレームにおける口形状を各音素に対応する口形状の重み付け合成として表現する。ここで、1401~1403は各フレームにおける音素毎の重み配分を示している。また、図1404~1406は、各々の重み付け合成によって生成される口形状を示している。

【0030】また、図15に示すように、子音発音時において、前後の口形状(1401, 1503)から、音素間のアニメーション(1502)を作成する。このとき、破裂音など、音素が実際に発音される前に口形状が変化しているものに関しては、イベントを前に移動(1504)しておく。また、アニメーション時間として、移動後のフレーム時間を用いるとする。ここで、図1505は/p/の発音イベントを示しており、この時、アニメーション(1502)に割り振っておく。また、図1506は/a/の発音イベントを示している。

【0031】図16に示すように、図15において作成した音素アニメーションとイベントデータの発生(1601)に従って、口形状のアニメーションを起動していく(1602)。次に、移動中のアニメーションの合成(1603)したあと、画面への表示(1604)を行う。

【0032】なお、アニメーションの合成に関しては、例えば、複数の基本のアニメーションを所定の係数で線形加算して目的のアニメーションを合成することができる。線形加算による合成する手法は、周知の手法を採用できる。例えば、特開平2000-11199号公報「アニメーションの自動生成方法」を採用してもよい。

【0033】アニメーション表示の様子を図17に示

す。まず、音声データ（1713）の再生を行い、これにマッピングされているイベント（1710、1711、1712）とアニメーションを起動していく。次に、起動したアニメーションの合成を行い（1708）最終的に口形状のアニメーションを表示する（1701、1702、1703）。ここで、1709はイベントデータを示しており、また1705、1706、1707は各アニメーションを示している。

【0034】なお、ストリーミングによって音声データが配信される場合には、音素解析が行われている時間分バッファリング（遅れ）が生じるが、これは計算機の能力向上に伴って遅れ時間を短くすることが可能である。

【0035】

【発明の効果】前記のように、本発明によれば、音素データに対するリップシンクアニメーションを予め作成しておくのではなく、音声データから適宜、合成によって生成するために、後から音声データを作成し、追加することが容易になる。また、音素解析におけるあいまいさに対してあいまいな口表現を行うという意味を付けることによって、より自然な口の表現が可能になる。

【図面の簡単な説明】

【図1】 本発明の実施例の実現例を全体として示す図である。

【図2】 上述実施例の概要を説明する図である。

【図3】 上実施例の動作の流れの概要を示す図である。

【図4】 上述実施例の音素イベントの生成を説明する図である。

【図5】 上述実施例における音声データの入力ソースを説明する図である。

【図6】 上述実施例における音素解析用の解析フレー

ムを説明する図である。

【図7】 上述実施例の解析フレーム単位の周波数スペクトルの取得を説明する図である。

【図8】 周波数スペクトルのフォルマント位置のキャリブレーションを説明する図である。

【図9】 解析フレーム単位でフォルマントを取り出すことを説明する図である。

【図10】 音素間の状態遷移図を説明する図である。

【図11】 同じ口のアニメーションとして集約できる状態遷移を一つにまとめる態様を説明する図である。

【図12】 計測されたフォルマントと、キャリブレーションとの関係から重みを計算する態様を説明する図である。

【図13】 子音に関して重み計算を行う態様を説明する図である。

【図14】 音素に対応した口の形状の合成状態としてあいまいな口を表現することを示す図である。

【図15】 音素間における口形状のアニメーションを作成することを示す図である。

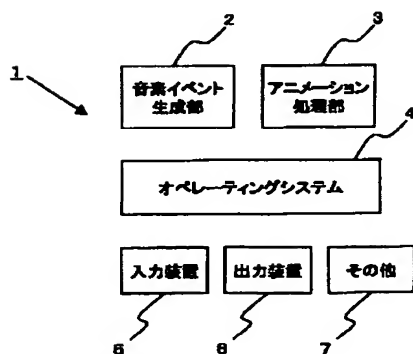
【図16】 くりアニメーションの合成と表示の流れを説明する図である。

【図17】 音声データの再生に合わせて口形状を表すアニメーションを合成していく様子を説明する図である。

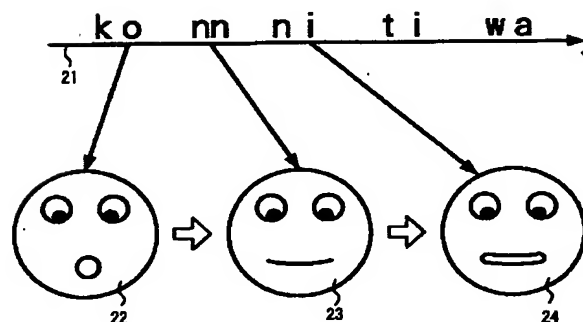
【符号の説明】

- 1 アニメーション生成装置
- 2 音素イベント生成部
- 3 アニメーション処理部
- 4 オペレーティングシステム
- 5 入力装置
- 6 出力装置

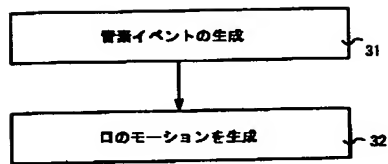
【図1】



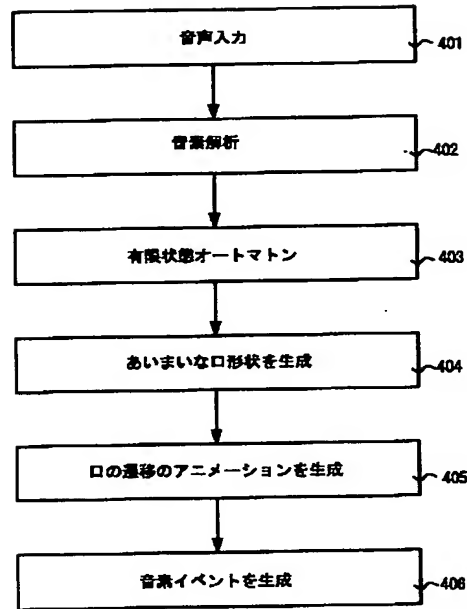
【図2】



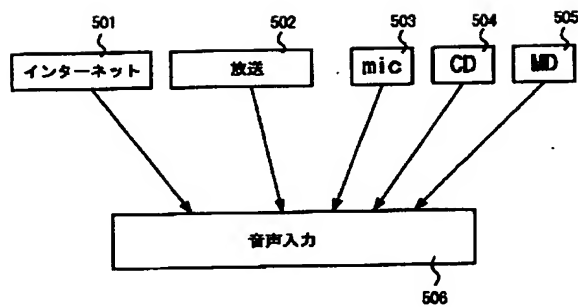
【図3】



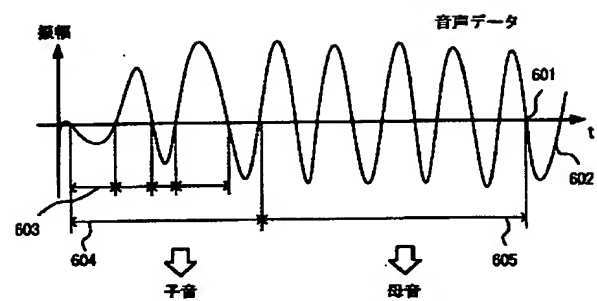
【図4】



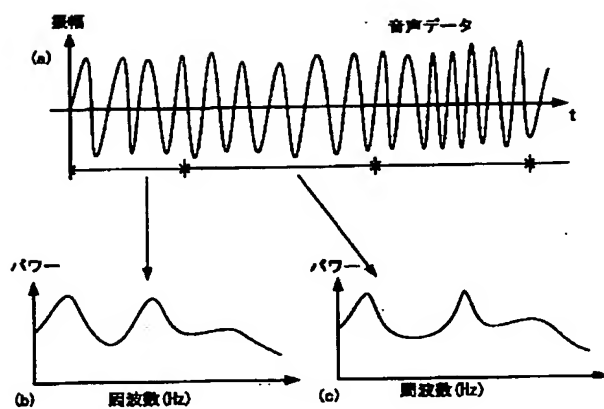
【図5】



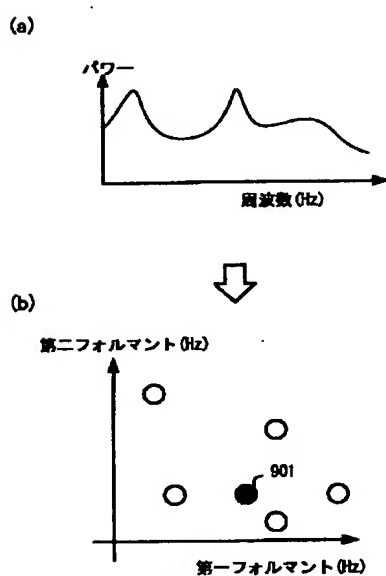
【図6】



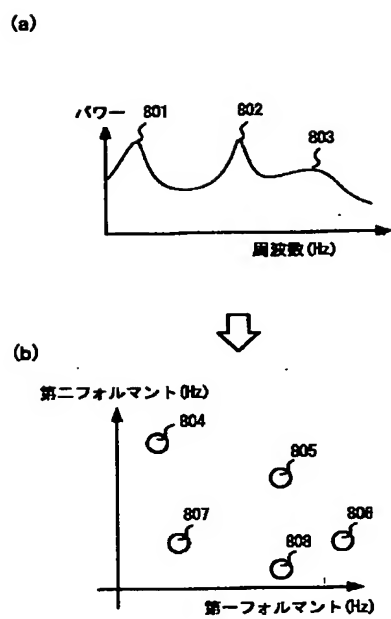
【図7】



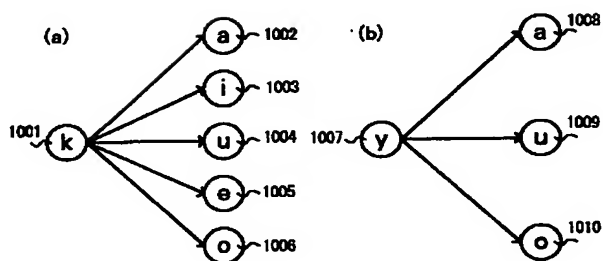
【図9】



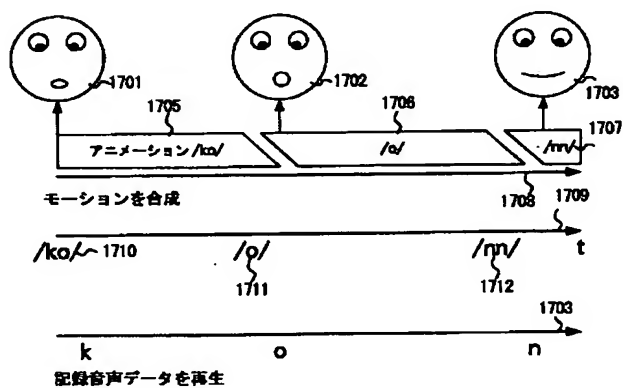
【図8】



【図10】

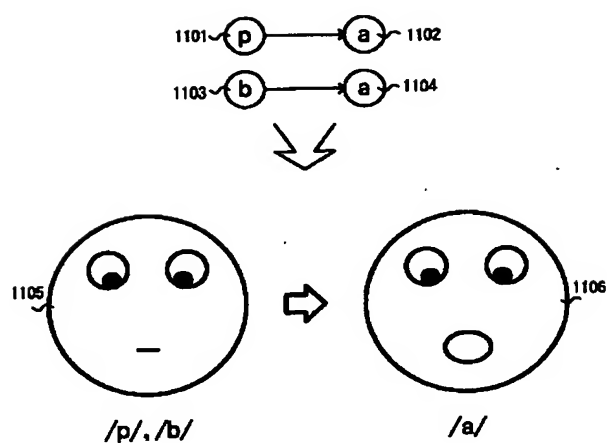


【図17】

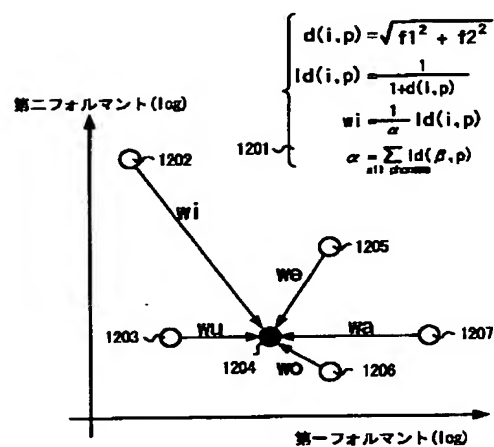




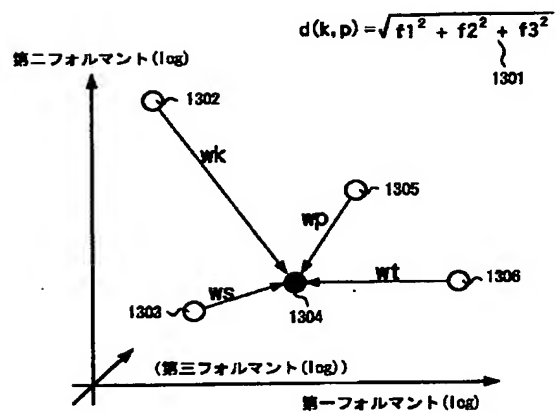
【図11】



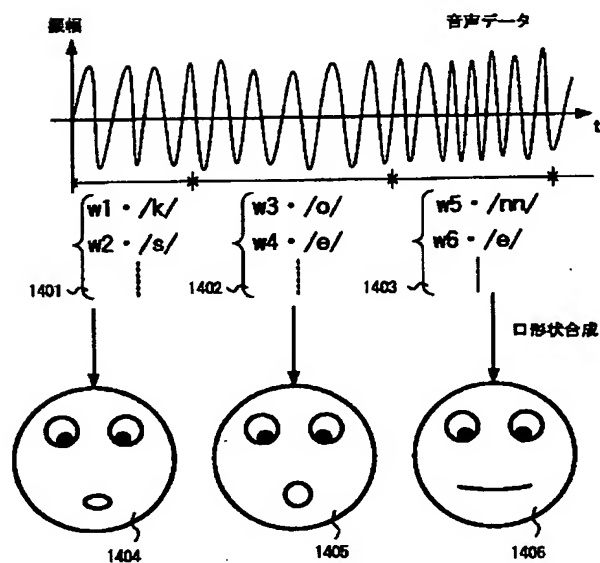
【図12】



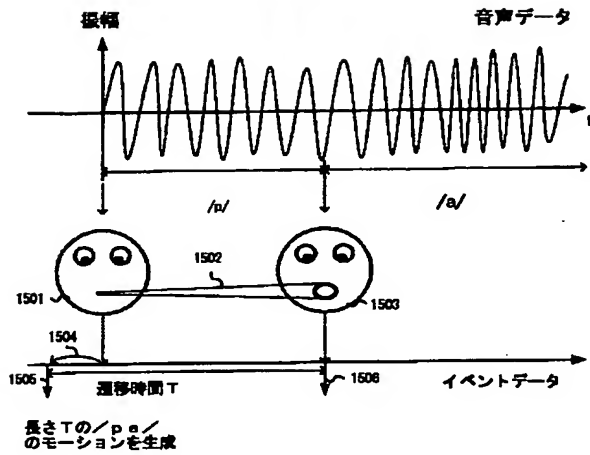
【図13】



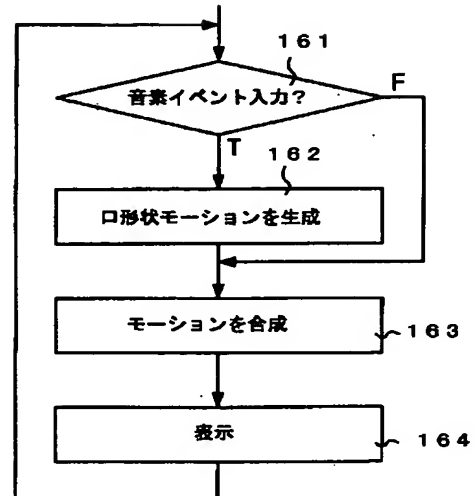
【図14】



【図15】



【図16】



フロントページの続き

(51) Int. Cl.<sup>7</sup>

識別記号

F I

テーマコード (参考)

G 1 0 L 21/06

G 1 0 L 3/00

S

15/22

5 6 1 C

F ターム (参考) 5B050 BA08 BA12 EA19 EA24 FA10

5D015 AA05 BB02 CC03 CC04 DD02

KK01 LL12

5D045 AB01 AB11